

**What Is Claimed Is:**

1. A method of flow controlling InfiniBand receive traffic,  
comprising:

- 5 maintaining a single memory structure for queuing InfiniBand traffic  
received via multiple virtual lanes and multiple queue pairs;  
identifying a first packet payload received via a first virtual lane and a first  
queue pair;  
determining whether the first payload can be stored in the memory  
10 structure without exceeding a portion of the memory structure allocated to the first  
virtual lane;  
determining whether the first payload can be stored in the memory  
structure without exceeding a portion of the memory structure allocated to the first  
queue pair; and  
15 if storing the first payload in the memory structure would exceed said  
portion of the memory structure allocated to the first queue pair, determining  
whether the first queue pair is enabled to use a shared portion of the memory  
structure to store payloads of packets received via the first queue pair.

- 20 2. The method of claim 1, further comprising:  
allocating a portion of the memory structure to each of the multiple virtual  
lanes; and  
allocating a portion of the memory structure to each of the multiple queue  
pairs.

- 25 3. The method of claim 1, wherein the memory structure comprises a  
set of linked lists of memory structure buffers, including one linked list for each of

the multiple queue pairs that are active.

4. The method of claim 1, further comprising:  
dropping the first payload if the first payload cannot be stored in the  
5 memory structure without exceeding the portion of the memory structure allocated  
to the first virtual lane.

5. The method of claim 1, further comprising:  
issuing a Retry, Not Ready, Negative Acknowledgement (RNR-NAK) if:  
10 the first payload cannot be stored in the memory structure without  
exceeding a portion of the memory structure allocated to the first queue  
pair; and  
the first queue pair is not enabled to use the shared portion of the  
memory structure.

15 6. The method of claim 1, further comprising:  
issuing a Retry, Not Ready, Negative Acknowledgement (RNR-NAK) if:  
the first payload cannot be stored in the memory structure without  
exceeding a portion of the memory structure allocated to the first queue  
20 pair;  
the first queue pair is enabled to use the shared portion of the  
memory structure; and  
the shared portion of the memory structure is full.

25 7. The method of claim 1, further comprising:  
defining one or more dedicated thresholds in the portion of the memory  
structure allocated to the first queue pair; and

for each of said dedicated thresholds, identifying a number of message credits the queue pair may advertise when the amount of the memory structure used by the queue pair exceeds said dedicated threshold.

5           8.       The method of claim 1, further comprising:  
              defining one or more shared thresholds in the shared portion of the  
memory structure; and  
              for each of said shared thresholds, identifying a number of message credits  
the queue pair may advertise when the amount of the shared portion used by the  
10       multiple queue pairs exceeds said shared threshold.

              9.       The method of claim 1, further comprising:  
              receiving a request on a second queue pair to perform an RDMA (Remote  
Direct Memory Access) Read operation; and  
15       based on an amount of data expected to be received via the RDMA Read  
operation, reserving a sufficient number of buffers in the memory structure.

              10.      The method of claim 1, further comprising:  
              in the single memory structure, reassembling the queued InfiniBand traffic  
20       into outbound communications;  
              receiving a payload on an idle queue pair, wherein a queue pair is idle if no  
traffic from the queue pair is stored in the single memory structure; and  
              only queuing the payload in the single memory structure if sufficient space  
in the single memory structure is reserved for completing reassembly of outbound  
25       communications on each non-idle queue pair.

              11.      A computer readable medium storing instructions that, when

executed by a computer, cause the computer to perform a method of flow  
controlling InfiniBand receive traffic, the method comprising:

maintaining a single memory structure for queuing InfiniBand traffic  
received via multiple virtual lanes and multiple queue pairs;

5 identifying a first packet payload received via a first virtual lane and a first  
queue pair;

determining whether the first payload can be stored in the memory  
structure without exceeding a portion of the memory structure allocated to the first  
virtual lane;

10 determining whether the first payload can be stored in the memory  
structure without exceeding a portion of the memory structure allocated to the first  
queue pair; and

if storing the first payload in the memory structure would exceed said  
portion of the memory structure allocated to the first queue pair, determining  
15 whether the first queue pair is enabled to use a shared portion of the memory  
structure to store payloads of packets received via the first queue pair.

12. The computer readable medium of claim 11, wherein the method  
further comprises:

20 defining one or more dedicated thresholds in the portion of the memory  
structure allocated to the first queue pair; and

for each of said dedicated thresholds, identifying a number of message  
credits the queue pair may advertise when the amount of the memory structure  
used by the queue pair exceeds said dedicated threshold.

25

13. The computer readable medium of claim 11, wherein the method  
further comprises:

defining one or more shared thresholds in the shared portion of the memory structure; and

for each of said shared thresholds, identifying a number of message credits the queue pair may advertise when the amount of the shared portion used by the multiple queue pairs exceeds said shared threshold.

14. The computer readable medium of claim 11, wherein the method further comprises issuing a Retry, Not Ready, Negative Acknowledgement (RNR-NAK) only if one of:

10 (a) the first payload cannot be stored in the memory structure without exceeding a portion of the memory structure allocated to the first queue pair; and

the first queue pair is not enabled to use the shared portion of the memory structure; and

15 (b) the first payload cannot be stored in the memory structure without exceeding a portion of the memory structure allocated to the first queue pair;

the first queue pair is enabled to use the shared portion of the memory structure; and

20 the shared portion of the memory structure is full.

15. A method of flow controlling InfiniBand traffic received via multiple queue pairs, the method comprising:

25 maintaining a single memory structure for queuing InfiniBand traffic received via multiple queue pairs;

for each of the multiple queue pairs, dedicating zero or more buffers in the memory structure for storing InfiniBand traffic received via said queue pair;

for each queue pair having a number of dedicated buffers greater than  
zero:

identifying one or more buffer thresholds between 0 and N,  
inclusive, wherein N is the number of buffers dedicated to the queue pair;  
5 and

for each of said buffer thresholds, identifying an amount of credits  
said queue pair may advertise after consuming said threshold level of  
buffers; and

within a shared portion of the single memory structure usable by a subset  
10 of the multiple queue pairs:

identifying one or more shared thresholds; and

for each of said shared thresholds, identifying an amount of credits  
each queue pair in said subset of queue pairs may advertise after said  
shared threshold level of buffers is consumed.

15

16. The method of claim 15, further comprising:

receiving a payload of a first packet, from a first queue pair, for storage in  
the memory structure; and

determining if the payload can be stored without exceeding said number of  
20 buffers dedicated to the first queue pair.

17. The method of claim 16, further comprising:

issuing a Retry, Not Ready, Negative Acknowledgement (RNR-NAK).

25

18. The method of claim 16, further comprising:

if the payload cannot be stored without exceeding said number of buffers  
dedicated to the first queue pair, determining whether the first queue pair is

permitted to use said shared portion of the memory structure; and

if the first queue pair is permitted to use said shared portion of the memory structure, determining whether the payload can be stored without exceeding a final threshold in said shared portion.

5

19. The method of claim 18, further comprising:

issuing a Retry, Not Ready, Negative Acknowledgement (RNR-NAK).

20. The method of claim 15, further comprising:

10 receiving an instruction to perform an RDMA (Remote Direct Memory Access) Read operation on a first queue pair; and

reserving a sufficient numbers of buffers in the memory structure to store data to be retrieved via the RDMA Read operation.

15 21. The method of claim 20, further comprising, before said reserving:

calculating the sum of the number of buffers in the memory structure currently used by the first queue pair and the number of buffers in the memory structure currently reserved for other RDMA Read operations for the first queue pair; and

20 comparing said sum to said number of buffers dedicated to the first queue pair.

22. The method of claim 21, further comprising:

25 determining whether the first queue pair is permitted to use said shared portion of the memory structure.

23. A computer readable medium storing instructions that, when

executed by a computer, cause the computer to perform a method of flow controlling InfiniBand traffic received via multiple queue pairs, the method comprising:

- 5 maintaining a single memory structure for queuing InfiniBand traffic received via multiple queue pairs;
  - for each of the multiple queue pairs, dedicating zero or more buffers in the memory structure for storing InfiniBand traffic received via said queue pair;
  - for each queue pair having a number of dedicated buffers greater than zero:
    - 10 identifying one or more buffer thresholds between 0 and N, inclusive, wherein N is the number of buffers dedicated to the queue pair; and
    - for each of said buffer thresholds, identifying an amount of credits said queue pair may advertise after consuming said threshold level of buffers; and
    - 15 within a shared portion of the single memory structure usable by a subset of the multiple queue pairs:
      - identifying one or more shared thresholds; and
      - for each of said shared thresholds, identifying an amount of credits
      - 20 each queue pair in said subset of queue pairs may advertise after said shared threshold level of buffers is consumed.

24. The computer readable medium of claim 23, wherein the method further comprises:

- 25 receiving an instruction to perform an RDMA (Remote Direct Memory Access) Read operation on a first queue pair; and
- reserving a sufficient numbers of buffers in the memory structure to store



data to be retrieved via the RDMA Read operation.

25. The computer readable medium of claim 23, wherein the method further comprises:

- 5 receiving a payload of a first packet, from a first queue pair, for storage in the memory structure;
- determining if the payload can be stored without exceeding said number of buffers dedicated to the first queue pair;
- if the payload cannot be stored without exceeding said number of buffers
- 10 dedicated to the first queue pair, determining whether the first queue pair is permitted to use said shared portion of the memory structure; and
- if the first queue pair is permitted to use said shared portion of the memory structure, determining whether the payload can be stored without exceeding a final threshold in said shared portion.

15

26. A method of flow controlling InfiniBand traffic received via multiple virtual lanes, comprising:

- maintaining a single memory structure for queuing InfiniBand traffic received via multiple active virtual lanes;
- 20 for each of the multiple active virtual lanes, identifying a threshold number of buffers in the memory structure dedicated to storing InfiniBand traffic received via said virtual lane; and
- for each of the multiple active virtual lanes:
  - receiving a value representing the total blocks sent on the virtual
  - 25 lane from the link partner of said virtual lane;
  - storing said value;
  - incrementing said value for each payload received via said virtual

lane and stored in the memory structure;

adding to said value the number of buffers in said dedicated set of buffers that are unused to calculate a credit limit; and advertising said credit limit to the link partner.

5

27. The method of claim 26, further comprising:

for each of the multiple active virtual lanes, maintaining a used buffer count to track the number of memory structure buffers used to store packet payloads received via said virtual lane.

10

28. The method of claim 27, further comprising:

receiving a payload of a packet received on a first virtual lane; determining whether the payload can be stored in the memory structure without causing the used buffer count to exceed said threshold number of buffers.

15

29. A computer readable medium storing instructions that, when executed by a computer, cause the computer to perform a method of flow controlling InfiniBand traffic received via multiple virtual lanes, the method comprising:

20

maintaining a single memory structure for queuing InfiniBand traffic received via multiple active virtual lanes;

for each of the multiple active virtual lanes, identifying a threshold number of buffers in the memory structure dedicated to storing InfiniBand traffic received via said virtual lane; and

25

for each of the multiple active virtual lanes:

receiving a value representing the total blocks sent on the virtual lane from the link partner of said virtual lane;

storing said value;  
incrementing said value for each payload received via said virtual  
lane and stored in the memory structure;  
adding to said value the number of buffers in said dedicated set of  
5 buffers that are unused to calculate a credit limit; and  
advertising said credit limit to the link partner.

30. A method of avoiding locking, in receive InfiniBand queues, the  
method comprising:  
10 maintaining a single memory structure for reassembling InfiniBand traffic  
received via multiple virtual lanes and multiple queue pairs;  
identifying a first packet payload received via a first queue pair that is idle,  
wherein the first queue pair is considered idle if no traffic from the first queue pair  
is stored in said single memory structure;  
15 for each other queue pair for which traffic from said queue pair is stored in  
said single memory structure, determining whether sufficient space in the single  
memory structure is reserved for reassembling said traffic; and  
storing the first packet payload in said single memory structure only if  
sufficient space in the single memory structure is available for reassembling said  
20 traffic.

31. The method of claim 30, wherein said determining comprises, for  
each said other queue pair:  
identifying an amount of space in said single memory structure reserved  
25 for said other queue pair; and  
comparing said amount of reserved space to an amount of space expected  
to be needed to complete reassembly of said traffic from said other queue pair.

32. An apparatus for flow controlling received InfiniBand traffic,  
comprising:

5 a single memory structure configured to queue payloads of InfiniBand  
traffic received via multiple virtual lanes and multiple queue pairs;  
a resource manager configured to manage the memory structure;  
a first module configured to facilitate the advertisement of virtual lane  
credits; and  
10 a second module configured to facilitate the advertisement of queue pair  
credits.

33. The apparatus of claim 32, wherein said single memory structure  
comprises multiple linked lists of memory structure buffers, including one linked  
list for each of the multiple queue pairs that is active.  
15

34. The apparatus of claim 32, wherein said first module comprises an  
InfiniBand link core.

35. The apparatus of claim 32, wherein said second module comprises  
20 an acknowledgement generator configured to generate transport layer  
acknowledgements.

36. The apparatus of claim 32, further comprising a processor interface  
configured to facilitate the programming of operating parameters associated with  
25 the multiple virtual lanes and the multiple queue pairs.

37. The apparatus of claim 32, further comprising:

a first memory configured to store one or more parameters associated with operation of a first virtual lane.

38. The apparatus of claim 37, wherein said one or more parameters  
5 include:

a count of the number of memory structure buffers currently used to store payloads of packets received via the first virtual lane; and

a threshold, wherein a first packet is dropped if storing the payload of the first packet would cause said count to exceed said threshold.

10

39. The apparatus of claim 32, further comprising:

a second memory configured to store, for each of the multiple queue pairs that is active, one or more parameters associated with operation of said queue pair.

15

40. The apparatus of claim 39, wherein said one or more parameters include:

a maximum number of message credits advertisable by said queue pair;

a maximum number of memory structure buffers dedicated to storing  
20 payloads of packets received via said queue pair; and

an indicator configured to indicate whether said queue pair is enabled to use a set of shared memory structure buffers.

41. The apparatus of claim 40, wherein said one or more parameters  
25 further include:

one or more dedicated thresholds, wherein each said dedicated threshold identifies a subset of said maximum number of memory structure buffers; and

for each said dedicated threshold, a number of message credits  
advertisable by said queue pair when said queue pair uses said subset of said  
maximum number of memory structure buffers.

5           42.    The apparatus of claim 40, wherein said one or more parameters  
further include:

          a number of shared memory structure buffers in said set of shared memory  
structure buffers, wherein said shared memory structure buffers are available for  
use by said queue pair to store payloads of packets received via said queue pair if:

10               said queue pair has used said maximum number of memory  
structure buffers; and

          said indicator indicates that said queue pair is enabled to use said  
set of shared memory structure buffers; and

          a maximum number of message credits advertisable by said queue pair  
15   when said queue pair starts using said shared memory structure buffers.

          43.    The apparatus of claim 42, wherein said one or more parameters  
further include:

          one or more shared thresholds, wherein each said shared threshold  
20   identifies a subset of said number of shared memory structure buffers; and

          for each said shared threshold, a number of message credits advertisable  
by said queue pair when said queue pair uses said subset of said number of shared  
memory structure buffers.

25           44.    An apparatus for flow controlling received InfiniBand traffic,  
comprising:

          a single memory structure comprising multiple linked lists of memory

structure buffers, including one linked list for each of multiple active queue pairs;  
and

a queue pair memory configured to store, for each of the multiple active queue pairs, a set of parameters for facilitating use of said single memory

5 structure, said parameters including:

a maximum number of message credits issuable for the queue pair;

a maximum number of memory structure buffers dedicated to  
storing payloads of packets received via the queue pair;

10 one or more thresholds less than said maximum number of  
dedicated memory structure buffers;

for each of said one or more thresholds, a number of message  
credits issuable for the queue pair when said threshold number of memory  
structure buffers are in use; and

15 an indicator configured to indicate whether a set of shared memory  
structure buffers may be used after said maximum number of dedicated  
memory structure buffers are in use.

45. The apparatus of claim 44, further comprising:

a set of global parameters for facilitating use of said set of shared memory

20 structure buffers, the global parameters including:

a number of shared memory structure buffers in said set of shared  
memory structure buffers;

a maximum number of message credits issuable by a queue pair  
using said set of shared memory structure buffers;

25 one or more shared thresholds less than said number of shared  
memory structure buffers; and

for each said shared threshold, a number of message credits

issuable by a queue pair using said set of shared memory structure buffers when the number of shared memory structure buffers in use exceeds said shared threshold.

5

46. The apparatus of claim 44, further comprising:

a resource manager configured to manage use of said single memory structure;

a receive packet processor configured to request one or more memory structure buffers from said resource manager to store a payload of a newly

10

received packet; and

a post packet processor configured to request reservation of one or more memory structure buffers for data to be received via an RDMA (Remote Direct Memory Access) Read operation.

15

47. The apparatus of claim 46, wherein said resource manager comprises said queue pair memory.

48. The apparatus of claim 44, further comprising:

20

a virtual lane memory configured to store a set of parameters for facilitating queuing, in said single memory structure, of payloads of packets received via multiple virtual lanes.

25

49. A method of preventing memory locks when storing traffic from multiple communication streams in a single memory structure, the method comprising:

for each of a set of communication streams received at a communication interface, maintaining in a single memory structure a queue for assembling



payloads of packets received via the communication streams into outbound communications;

receiving a first packet via a first communication stream, wherein a payload of the packet comprises a portion of a first outbound communication to

5 be assembled;

determining a combined depth of the queues;

determining if a first threshold amount of space in the single memory structure is free; and

if less than said first threshold amount of space is free:

10 determining if another portion of the first outbound communication is stored in the single memory structure; and

if no other portion of the first outbound communication is stored in the single memory structure, rejecting the first packet.

15 50. The method of claim 49, wherein said first threshold amount of space is approximately equal to a Maximum Transfer Unit (MTU) in effect for one or more communication streams in the set of communication streams.

20 51. The method of claim 49, wherein said first threshold amount of space is approximately equal to the size of the first outbound communication.

52. The method of claim 49, further comprising, if less than said threshold amount of space is free:

25 identifying a subset of the multiple communication streams, wherein for each communication stream in the subset of communication streams an outbound communication is being assembled in the single memory structure, including a second outbound communication for a second communication stream; and

after the second outbound communication is assembled, preventing storage of a portion of another outbound communication for the second communication stream until outbound communications for a threshold number of the subset of communication streams are assembled.

5

53. The method of claim 49, further comprising if no other portion of the first outbound communication is stored in the single memory structure: recording rejection of the first packet.

10

54. The method of claim 49, further comprising: identifying communication streams for which a packet is rejected.

55. The method of claim 54, further comprising: awarding priority to the identified communication streams, for storing packets in the single memory structure, when said first threshold amount of space becomes free.

15